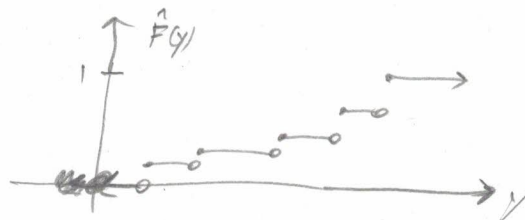


lec 07 MATH 383

This method will work for all parametric models. But what if you don't know the model? non-parametric (i.e. no model is assumed)

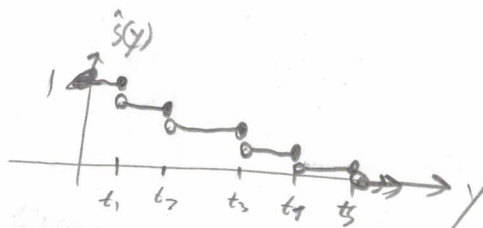
Let's go back to 381. Recall \hat{F} , the empirical CDF estimator for Y_1, \dots, Y_n iid.

$$\hat{F}(y) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i \leq y}$$



Thus \hat{S} , the empirical survival estimator is just its complement

$$\begin{aligned} \hat{S}(y) &:= 1 - \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i \leq y} \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i > y} \end{aligned}$$



Let the t_i 's be the order statistics of the Y_i 's. Assume all unique.

Thus, $0 < t_1 < t_2 < \dots < t_n$. Let's tabulate:

time	# still surviving (n_i)	# "at risk" at time of first death event
$t_0 = 0$	$n_0 = n$	$0 = n_{n+1} < n_n < \dots < n_2 < n_1 = n$
t_1	$n_1 = n-1$	
t_2	$n_2 = n-2$	
\vdots		
t_n	$n_n = n-n = 0$	

$$\hat{S}(0) = 1$$

$$\hat{S}(t_1) = \hat{P}(T > t_1) = \frac{\sum \mathbb{1}_{Y_i > t_1}}{n} = \frac{n_2}{n_1} = \frac{n-1}{n} = 1 - \frac{1}{n}$$

$$\hat{S}(t_2) = \hat{P}(T > t_2) = \hat{P}(T > t_2 | T > t_1) \hat{P}(T > t_1) = \frac{\sum \mathbb{1}_{Y_i > t_2}}{\sum \mathbb{1}_{Y_i > t_1}} \cdot \frac{n_2}{n_1} = \frac{n_3}{n_2} \cdot \frac{n_2}{n_1} = \frac{n_3}{n_1} = \frac{1}{n} \sum \mathbb{1}_{Y_i > t_2}$$

$$\begin{aligned} \hat{S}(t_3) &= \hat{P}(T > t_3) = \hat{P}(T > t_3 | T > t_2) \hat{P}(T > t_2 | T > t_1) \hat{P}(T > t_1) \\ &= \frac{n_4}{n_3} \cdot \frac{n_3}{n_2} \cdot \frac{n_2}{n_1} = \frac{n_4}{n_1} = \frac{1}{n} \sum \mathbb{1}_{Y_i > t_3} \\ &= \frac{n_3-1}{n_3} \cdot \frac{n_2-1}{n_2} \cdot \frac{n_1-1}{n_1} = \left(1 - \frac{1}{n_3}\right) \left(1 - \frac{1}{n_2}\right) \left(1 - \frac{1}{n_1}\right) \\ &= \left(1 - \frac{1}{n-2}\right) \left(1 - \frac{1}{n-1}\right) \left(1 - \frac{1}{n}\right) \end{aligned}$$

$$\hat{S}(t_k) = \prod_{i=k}^n \left(1 - \frac{1}{n_i}\right) = \prod_{i=1}^k \left(1 - \frac{1}{n-(i-1)}\right) \quad \text{"product limit estimator"}$$

If $j, k \in \{0, 1, \dots, n\}$

If $j < k$, $y \in (t_j, t_k) \Rightarrow \hat{S}(y) = \prod_{\{i: y < t_i\}} (1 - \frac{1}{n_i})$

How about inference for a certain value of y e.g. 1 yr?
 $\theta = P(Y > 1)$, probab. of surviving to 1 yr or more.

$$\hat{\theta} = \hat{S}(1) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i > 1} = \prod_{\{i: 1 < t_i\}} (1 - \frac{1}{n_i})$$

We've seen this before...

Let $B_i = \mathbb{1}_{Y_i > 1} \Rightarrow B_1, \dots, B_n \stackrel{iid}{\sim} \text{Bern}(E[\mathbb{1}_{Y_i > 1}]) = \text{Bern}(P(Y_i > 1)) = \text{Bern}(\theta)$

$$\Rightarrow \hat{\theta} = \bar{B} \underset{\substack{\uparrow \\ \text{by CLT}}}{\sim} N(E[B_i], \frac{\text{Var}(B_i)}{n}) = N(\theta, \frac{\theta(1-\theta)}{n})$$

By Slutsky's

$$\Rightarrow \hat{\theta} \sim N(\theta, \frac{30(1-30)}{n})$$

With this result, you can plot \pm bands around $\hat{S}(t)$

How about inference for a quantile? Eg. Median

$$\hat{\theta} = \hat{\text{Med}}(Y) = \argmin_y \{ \hat{S}(y) \leq \frac{1}{2} \}$$

$$\text{Var}(\hat{\theta})?$$

This is worked out asymptotically.

but beyond scope of course why?

Ans: $\hat{\theta} \sim N(\text{Med}(Y), \text{Var}(\hat{\theta}))$

For the purpose of this class, use the bootstrap!

Inference for the mean? Use $\bar{Y} \sim N(\theta, \frac{\sigma^2}{n})$

Compare 2-samples

$Y_{1,1}, \dots, Y_{1,n_1} \stackrel{iid}{\sim} \text{Dist}_1, Y_{2,1}, \dots, Y_{2,n_2} \stackrel{iid}{\sim} \text{Dist}_2$

Both OBP's
Same OBP's

Ans: OBP₁ \neq OBP₂

\Rightarrow Use 2-sample K-S test

What if there is non-unique survival times?

No change! However, there is an updated formula for the product limit estimator

Let $n' = \# \text{ of unique } t_0$

where $n' \leq n$ and $n' < n$ if

there are
least one duplicate

3

time	# of censored (d_i)	# still surviving (h_i)
$t_0 = 0$	0	$h_1 = n$
t_1	d_1	$h_2 = h_1 - d_1$
t_2	d_2	$h_3 = h_2 - d_2$
t_3	d_3	$h_4 = h_3 - d_3$
\vdots		
t_n	d_n	$h_{n+1} = h_n - d_n = 0$

$$0 = h_{n+1} < h_n < \dots < h_2 < h_1 = n$$

Note: the last quantity never enters into the equation

$$\hat{S}(t_3) = \hat{P}(T > t_3 | T > t_2) \hat{P}(T > t_2 | T > t_1) \hat{P}(T > t_1)$$

$$= \frac{h_4}{h_3} \cdot \frac{h_3}{h_2} \cdot \frac{h_2}{h_1} = \frac{h_4}{h_1} = \frac{1}{n} \sum \mathbb{1}_{y_i > t_3}$$

$$= \frac{h_3 - d_3}{h_3} \cdot \frac{h_2 - d_2}{h_2} \cdot \frac{h_1 - d_1}{h_1}$$

$$= \left(1 - \frac{d_3}{h_3}\right) \cdot \left(1 - \frac{d_2}{h_2}\right) \cdot \left(1 - \frac{d_1}{h_1}\right)$$

$$\hat{S}(t_k) = \prod_{i=0}^{k-1} \left(1 - \frac{d_i}{h_i}\right)$$

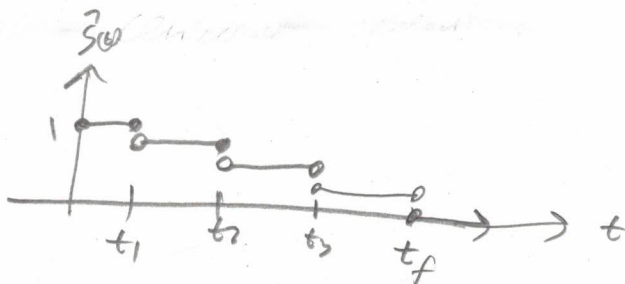
\downarrow # died in this time period
 \uparrow # "at risk" or still alive just before t_k

$$\hat{S}(t) = \prod_{\{i: y_i < t\}} \left(1 - \frac{d_i}{h_i}\right) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{y_i > t}$$

What about censoring?

If we have right-censoring at $t_f := \max(Y)$, then what happens?

$$\hat{S}(y) = \frac{1}{n} \mathbb{1}_{y_i > y}$$

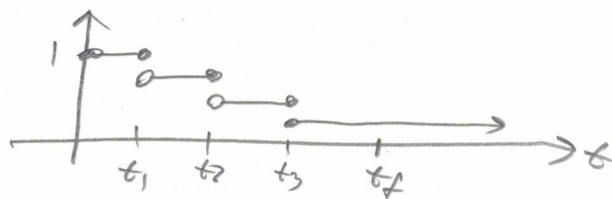


this is clearly wrong!

We know $\hat{S}(y) > 0$ for $y > t_f$. ^{some} censored observations are not the same as events!

Solution: just ignore the censoring!

$$\hat{S}(y) = \frac{1}{n} \mathbb{1}_{y_i > 0 \text{ and } c_i = 0}$$



Much more realistic.

Dropping all y_i s.t. $c_i = 1$ from the dataset changes nothing except...

\bar{y} is non biased downward you deleted lots of real y_i 's when $y_i > t_f \forall i$! No good way to get around this. This is why people focus on $O = \text{MED}(Y)$ and not $E(Y)$ since it's tractable. Others focus on $O = \text{restricted mean}$

$$:= \frac{\int_0^{t_f} y f(y) dy}{F(t_f)}$$
 which is a restricted survival $[0, t_f] \subset \mathcal{Y} = [0, \infty)$.
 but we won't study this